# Phylogenomic analysis of disease-resistance proteins in plants

Kimmen Sjölander
University of California Berkeley
http://phylogenomics.berkeley.edu

Collaborators:
David Konerding
Wayne Christopher
Bob Edgar
Joseph Dale
Austin Huang

Collaborators:
Brian Staskawicz
Barbara Baker
Richard Michelmore

# Challenges in protein classification

1. **Remote homolog detection.**
   *How much information does knowing a remote homolog provide?*

2. **Phylogenetic context is critical.**
   *Paralogs can have divergent function (so can orthologs...)*

3. **Domain structure issues.**

4. **Some fraction of the annotations in the sequence databases are not exactly accurate.**

# Function and Structure Prediction by Homology

If you have a sequence you know nothing about, find a relative.

# Given one member, find the relatives...

Would we recognize this member?



Homolog identification and profile construction helps differentiate critical features from variable
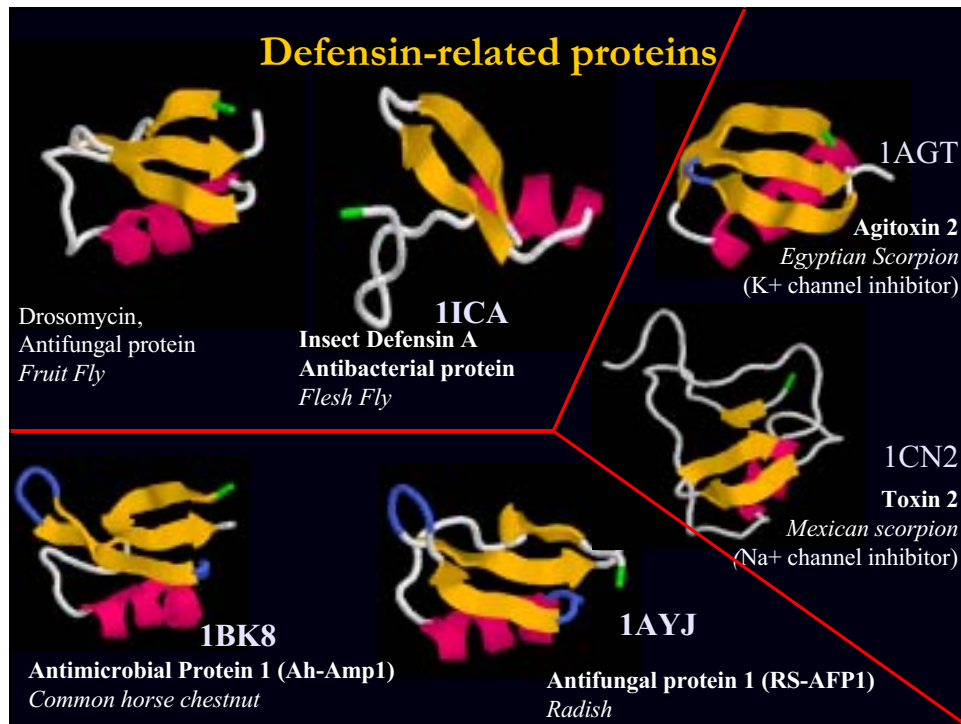
Profile generalization allows us to identify some truly remote relatives

**"Evolution conserves structure and function"**

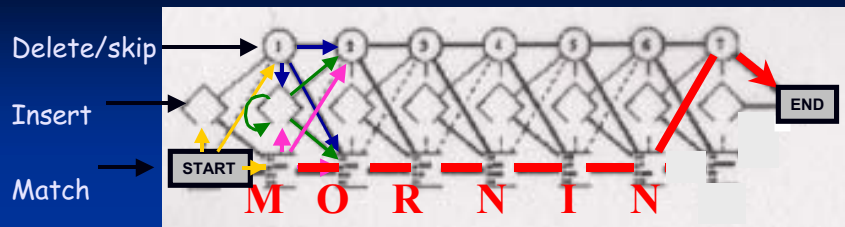But not completely.

Defensin-related proteins

---

# Why not just use BLAST?

Poor performance in remote homolog detection compared to HMM methods.

Park, J., Karplus, K., Barrett, C., Hughey, R., Haussler, D., Hubbard, T. & Chothia, C. Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods. J. Mol. Biol. 284, 1201-1210 (1998).

# Hidden Markov Model (HMM)



Delete/skip

Insert

Match

START    END

M O R N I N

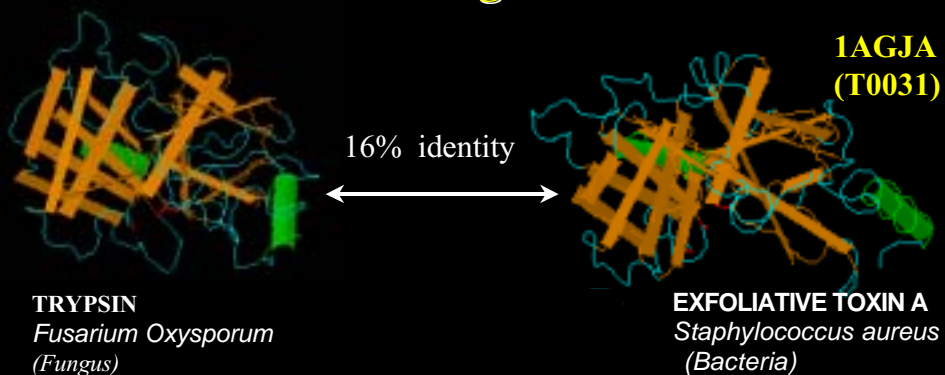Originally used in speech recognition (Rabiner, 1986)

Proposed for DNA modeling (Churchill, 1989)

Applied to modeling proteins (Haussler et al, 1992)

- Multiple sequence alignment
- Identification of related family members ("homologs")

Hidden Markov Models in Computational Biology: Applications to Protein Modeling. Krogh, Brown, Mian, Sjolander and Haussler, *J.Mol. Biol.* (1994)

---

# Homolog recognition in the Twilight Zone CASP2 Target T0031



1AGJA (T0031)

16%  identity

**TRYPSIN**
*Fusarium Oxysporum*
*(Fungus)*

**EXFOLIATIVE TOXIN A**
*Staphylococcus aureus*
*(Bacteria)*

For homolog recognition in the Twilight Zone, we need to know:
Which positions are critical?
Where can we allow deletions or mutations?

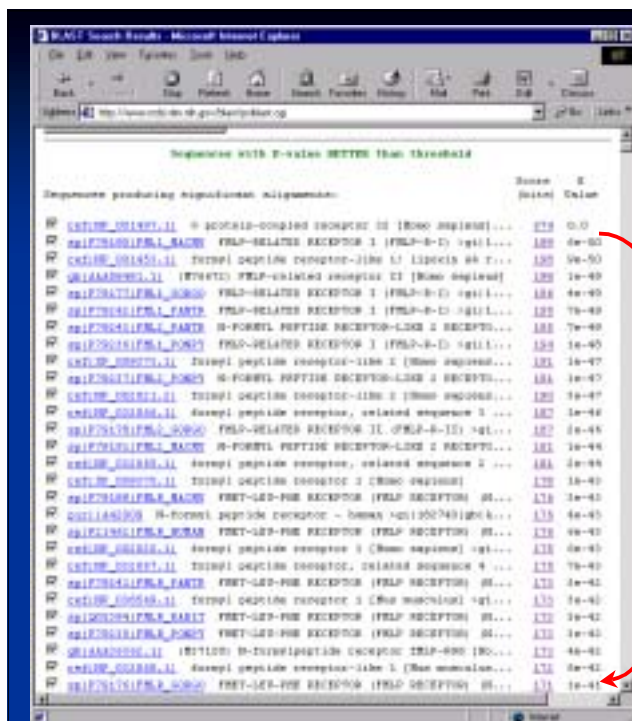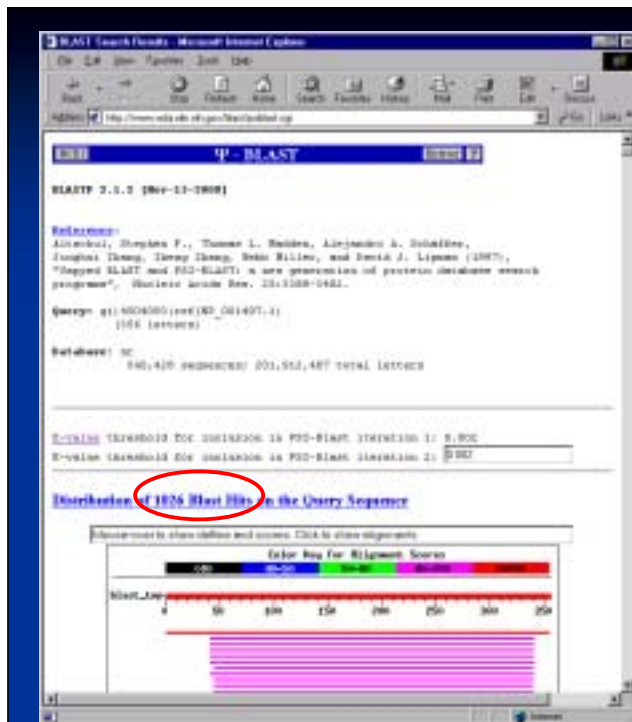# Homolog detection is just the first step…

Correct functional classification requires attention to evolutionary relationships

---

## Example 1: Orphan GPCR classification

Step 1:
Run
BLAST



Is the
query an
FMLP
receptor?

Or a chemokine receptor?

Or a C5A Anaphylatoxin chemotactic receptor?

Or…?



**Or a novel type of receptor?**

# Another reason to not rely on pairwise sequence similarity

## What if the top-scoring match is incorrectly annotated?

# Example 2: Errors in database annotations

# The top matching BLAST hits are also putative odorant receptors

# Lessons from CASP2



1. HMMs optimized for remote homolog detection generally require clustering and alignment of many divergent sequences.

2. Alignments of new sequences to these HMMs can be pretty awful.
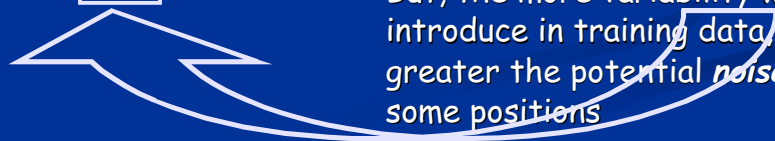
Given a protein sequence ("target"), predict its most likely fold, and produce an alignment of the target and the solved structure. Predictions judged using structure-structure alignment (SCOP, VAST, DALI).

---

# Conflict

```
D S L  F M  K I
D S I  F M  K V
D T I  W M  K M
D T I  W M  K L
D T V  W M  K F
D T F  R K  K I
D T F  R K  K V
```

- For effective *remote homolog* detection, a profile or HMM needs information from divergent family members

- Without this context, we cannot differentiate critical from variable positions

- HMMs constructed with such data provide a coarse classification

- But, the more variability we introduce in training data, the greater the potential *noise* at some positions

12

# Subfamily HMM construction

---

## How to build Subfamily HMMs (SHMMs)

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | D | S | L | F | M | K | I |
| 2 | D | S | I | F | M | K | V |
| 3 | D | T | I | W | M | K | M |
| 4 | D | T | I | W | M | K | L |
| 5 | D | T | V | W | M | K | F |
| 6 | D | T | F | R | K | K | I |
| 7 | D | T | F | R | K | K | V |

Share statistics between subfamilies where there is evidence of a common distribution..

Keep statistics separate at positions where there is evidence of *divergent* structure.

Improved specificity, sensitivity, alignment accuracy

# Step 1: Form Dirichlet Mixture Posterior

At each position, for each subfamily, construct a Dirichlet mixture *posterior*, by combining the Dirichlet mixture *prior* with the amino acids aligned at that position by that subfamily.

(Weighted) subfamily counts

Mixture coefficient

Component Parameters

$$q_j := P(\vec{\alpha}j \mid \vec{n}, \Theta^{Prior})$$

$$\alpha_{ji} := \alpha_{ji} + n_i$$

(Weighted) subfamily counts of amino acid *i*

# Step 2: Calculate family contribution

Other subfamilies contribute, proportional to the probability of the amino acids they aligned at that position, given the revised Dirichlet mixture density.

$$f_i = \sum_{s' \neq s} P(\vec{n}_{s'} \mid \Theta^{Post}) \, n_{s'i}$$

```
D S L F M K I
D S I F M K V
D T I W M K M
D T I W M K L
D T V W M K F
D T F R K K I
D T F R K K V
```

(Weighted) counts from subfamily s_

(Formula for computing Prob (n | _ ) are in Sjolander et al, 1996)

## Step 3: Compute pseudocounts

Add the family contribution to the observed (weighted) counts, to obtain the pseudocounts $t_i$ of amino acid $i$:

$$t_i = n_{si} + f_i$$

(Weighted) subfamily counts for subfamily s

family contribution
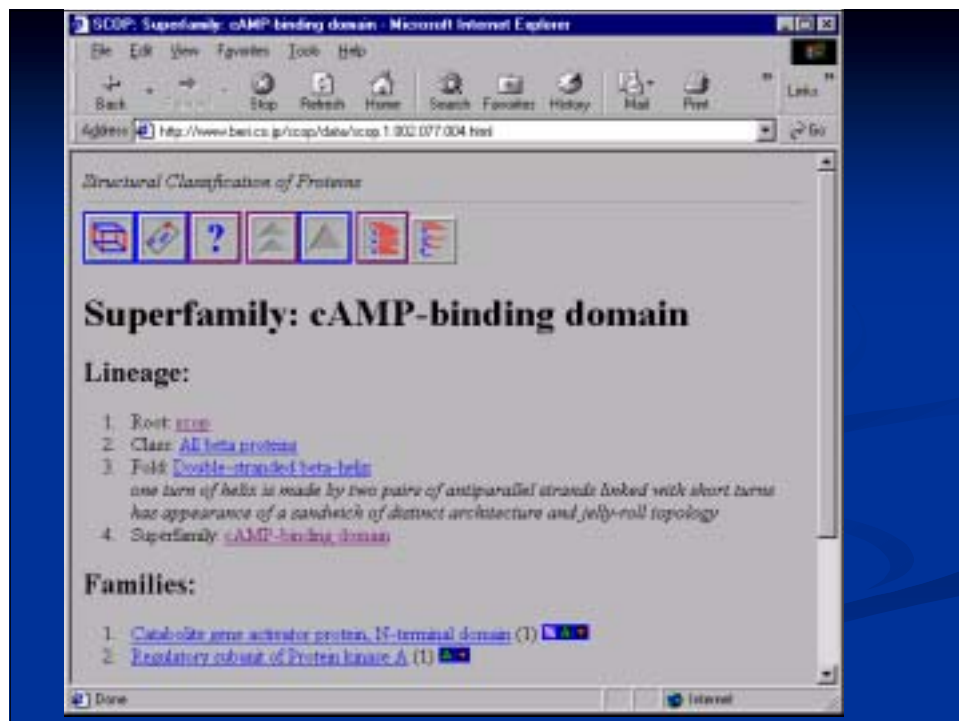
## Step 4: Compute amino acid probabilities

Normally, we compute amino acid probabilities by combining a Dirichlet mixture prior with observed counts as follows:

$$\hat{p}_i = \sum_j P(\bar{\alpha}_j \mid \vec{n}) \frac{n_i + \alpha_{ji}}{|\vec{n}| + |\bar{\alpha}_j|}$$

Instead, we will estimate the probability of amino acid $i$ as follows:

$$\hat{p}_i = \sum_j P(\bar{\alpha}_j \mid \vec{n}) \frac{t_i + \alpha_{ji}}{|\vec{t}| + |\bar{\alpha}_j|}$$

# Subfamily HMM Performance

# Socrates' First Command:
## Know Thyself

Test 1. How accurate are subfamily HMMs at recognizing their own training sequences?

---

# Test 1: Training Sequence Recognition

**Recognition of training sequences**



% sequences found above cutoffs: SHMM method  100%

GHMM method   99.89%

# Honor thy father and thy mother

*and thy brothers and sisters*
*and aunts and uncles*
*and cousins*
*and second cousins*
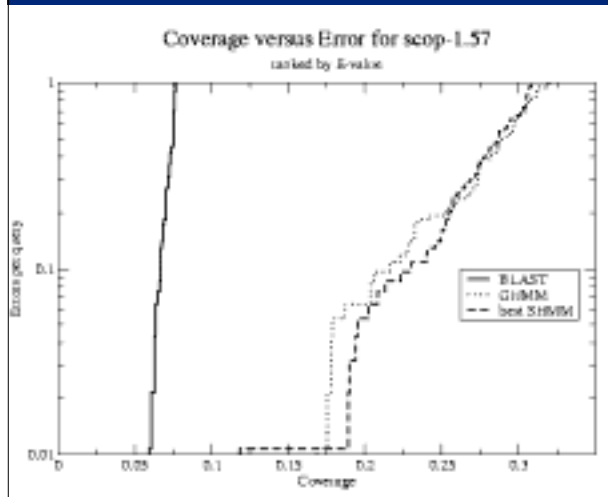*and third cousins twice removed…*

Recognition and classification of
family members

---

# Test 2: PSI-BLAST homolog detection
## Average Per Family

**Average Psi-BLAST homolog detection**



GHMM
SHMMs
ALL HMMs

Percentage scoring above cutoff

NLL-NULL score cutoff

## Subfamily HMMs improve homolog detection
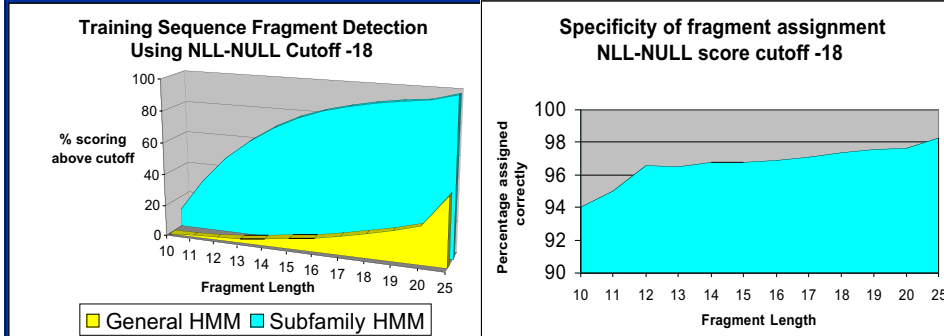### (relative to BLAST or a single HMM for the family)



PDB40 dataset 4,013 seqs
For each structure, homologs were gathered and aligned using SAM-T99 (from UCSC). A general HMM for each family was constructed from each alignment using Karplus sequence weighting and Dirichlet mixture densities. Subfamily HMMs were created from the same alignment. All PDB40 sequences were scored against each cluster, and assigned a general HMM score and the best Subfamily HMM (SHMM) score. Scores were sorted by significance. Homologs are determined by the SCOP database.

## Fragments and ESTs
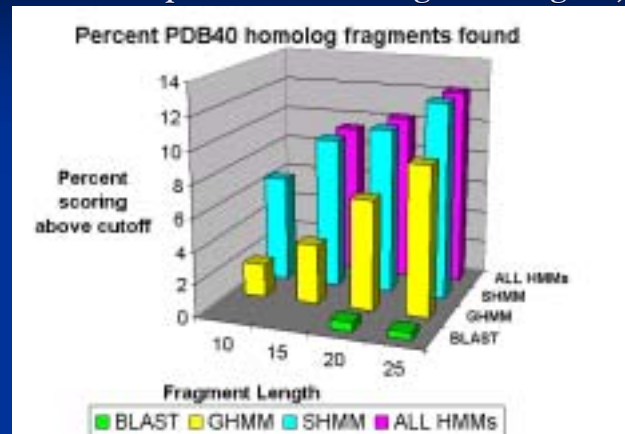## can be especially challenging

# Training sequence fragment detection

**Training Sequence Fragment Detection Using NLL-NULL Cutoff -18**



□ General HMM  □ Subfamily HMM

**Specificity of fragment assignment NLL-NULL score cutoff -18**



**PDB40 experiments.**
**Fixed cutoff chosen to provide zero FP for sequence lengths <= 25**

---

# PDB40 homolog fragment detection
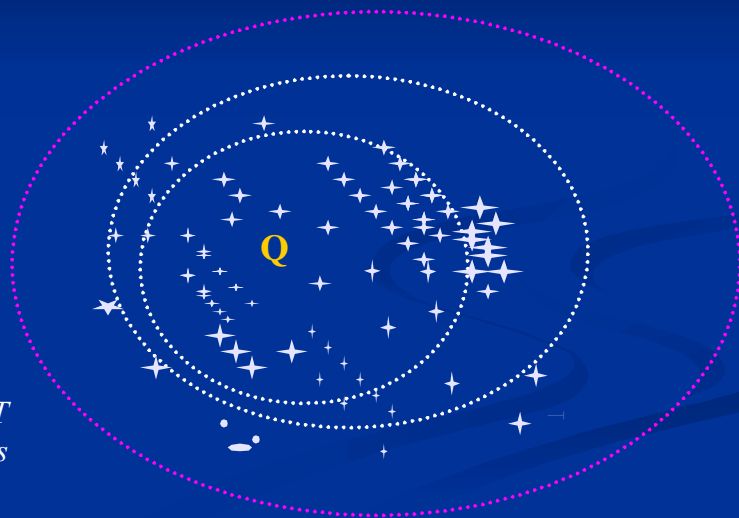### (HMM score and BLAST cutoffs chosen to give zero false positives for all fragment lengths )



| | 10 | 15 | 20 | 25 |
|---|---|---|---|---|
| Subfamily HMM / ALL: | 6.59% | 9.36% | 10.28% | 12.16% |
| General HMM: | 2.06 | 3.69 | 6.81 | 9.25 |
| BLAST: | | | 0.57 | 0.43 |

# FlowerPower

Iterative clustering and alignment tool

---

## Step 1: Identify putative homologs to query sequence (Q)

*PSI-BLAST*
*3 iterations*
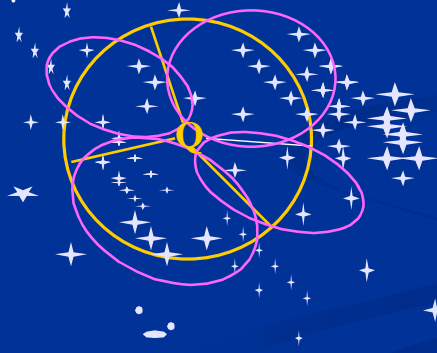*E-5 cutoff*

## Step 2: Select initial training set



## Step 3: Align initial set, identify subfamilies, and build subfamily HMMs.
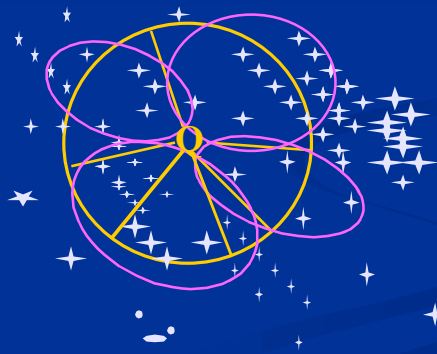
## Step 4: Identify and align new homologs.

1. Search with subfamily and general HMMs.
2. Accept hits above threshhold.
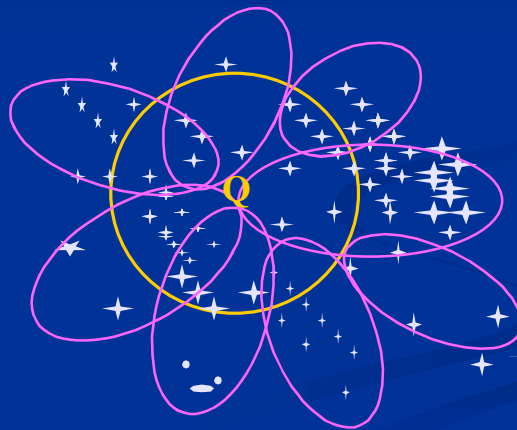3. Align accepted hits to closest HMM.



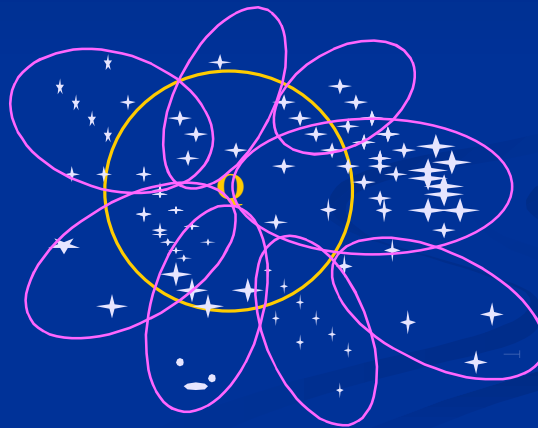## Step 5: Run BETE to identify subfamilies, and build new subfamily HMMs.
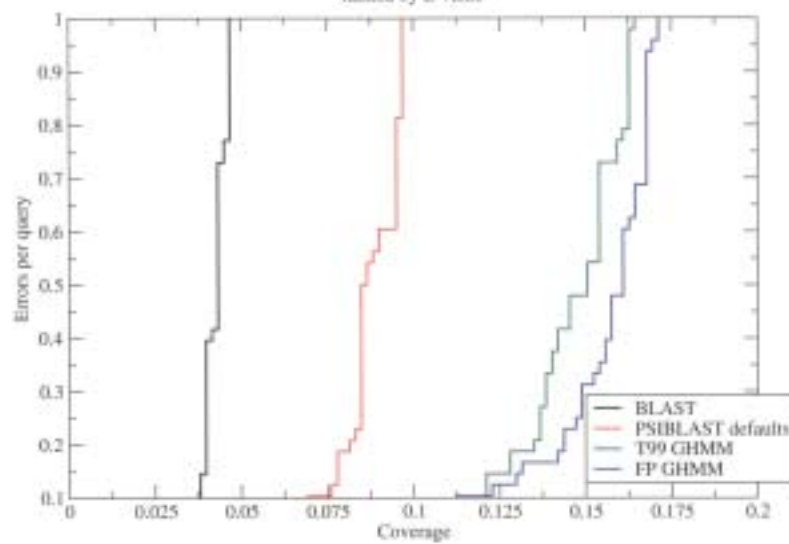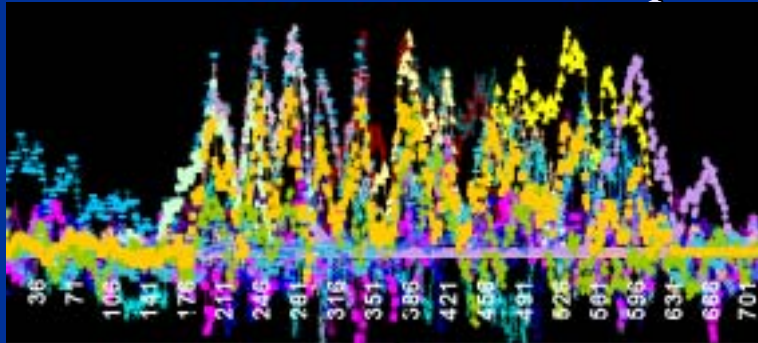
**Step 6: Iterate**



**until …**

convergence.

## HMM performance improves with improved alignments (Preliminary results)

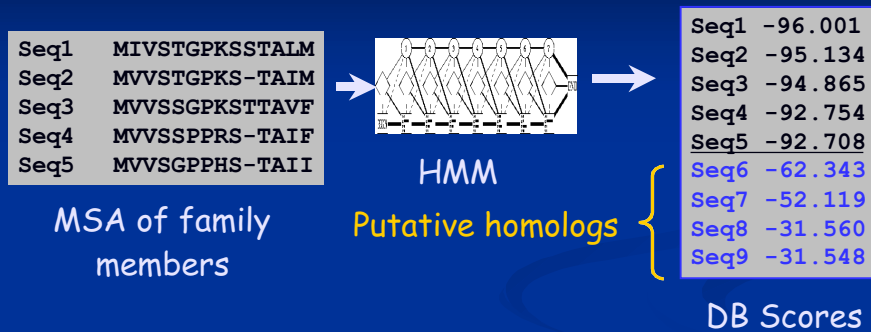Coverage versus Error for scop-1.57

ranked by E-value

Errors per query vs Coverage

Legend:
BLAST
PSIBLAST defaults
T99 GHMM
FP GHMM

# A Tale of Two Domains

**Hidden Markov models,
Potassium channels,
Excursions in the Twilight Zone
and some stories about startups…**

---

# Step 1: Build an HMM, and identify putative homologs

| Seq1 | MIVSTGPKSSTALM |
|------|----------------|
| Seq2 | MVVSTGPKS-TAIM |
| Seq3 | MVVSSGPKSTTAVF |
| Seq4 | MVVSSPPRS-TAIF |
| Seq5 | MVVSGPPHS-TAII |

MSA of family members

HMM

Putative homologs

```
Seq1 -96.001
Seq2 -95.134
Seq3 -94.865
Seq4 -92.754
Seq5 -92.708
Seq6 -62.343
Seq7 -52.119
Seq8 -31.560
Seq9 -31.548
```

DB Scores

Underlying assumption: Domain-level matches to the HMM will cluster in specific regions, while sequence fragments will align over HMM uniformly.

## Step 2: Align database hits to HMM

```
>Seq6
MIVSTSG
>Seq7
MVVTTG
>Seq8
SP
>Seq9
PP
```



| Seq6 | M | I | V | S | T | S | G |
|------|---|---|---|---|---|---|---|
| Seq7 | M | V | V | – | T | T | G |
| Seq8 | – | – | – | – | – | S | P |
| Seq9 | – | – | – | – | – | P | P |

(Alignment to subfamily HMMs can improve results)

## Step 3: Create Affinity (log odds) vectors

| Seq6 | M | I | V | S | T | S | G |
|------|---|---|---|---|---|---|---|
| Seq7 | M | V | V | – | T | T | G |
| Seq8 | – | – | – | – | – | S | P |
| Seq9 | – | – | – | – | – | P | P |

**Alignment of database hits**

| Seq6 | 3.3 | 3.2 | 3.3 | 3.4 | 3.5 | 1.5 | 2.8 |
|------|-----|-----|-----|-----|-----|-----|-----|
| Seq7 | 3.3 | 3.1 | 3.3 | 0.0 | 3.5 | 1.6 | 2.8 |
| Seq8 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.5 | 5.0 |
| Seq9 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 3.4 | 5.0 |

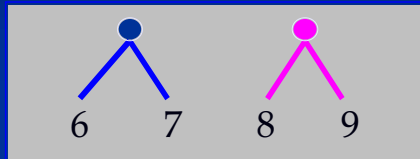**Log likelihood (Affinity) vectors**

$$L[s,p] = \log \frac{\text{Prob}(s_p \mid \text{HMM})}{\text{Prob}(s_p)}$$

$s_p$ = amino acid aligned by sequence $s$ at position $p$.
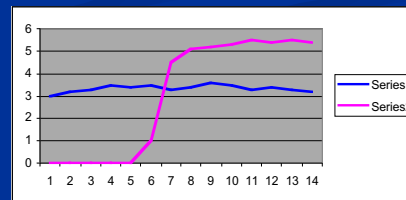Define $L[s,p] = 0$ when $s_p$ is a gap

| Seq1 | M | I | V | S | T | G | P |
|------|---|---|---|---|---|---|---|
| Seq2 | M | V | V | S | T | G | P |
| Seq3 | M | V | V | S | S | G | P |
| Seq4 | M | V | V | S | S | P | P |
| Seq5 | M | V | V | S | G | P | P |

**Alignment used as basis for HMM**

**Step 4: Cluster affinity vectors**
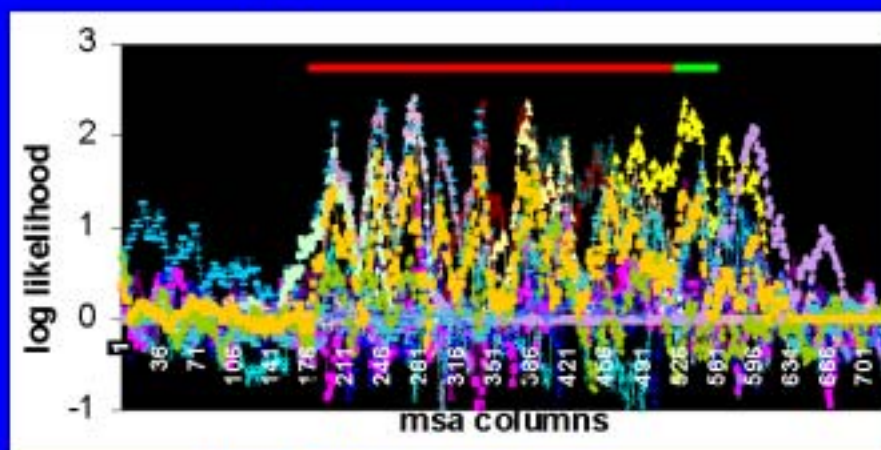
6  7    8  9

**Clustering of vectors**
**Agglomerative clustering**
**Euclidean distance**
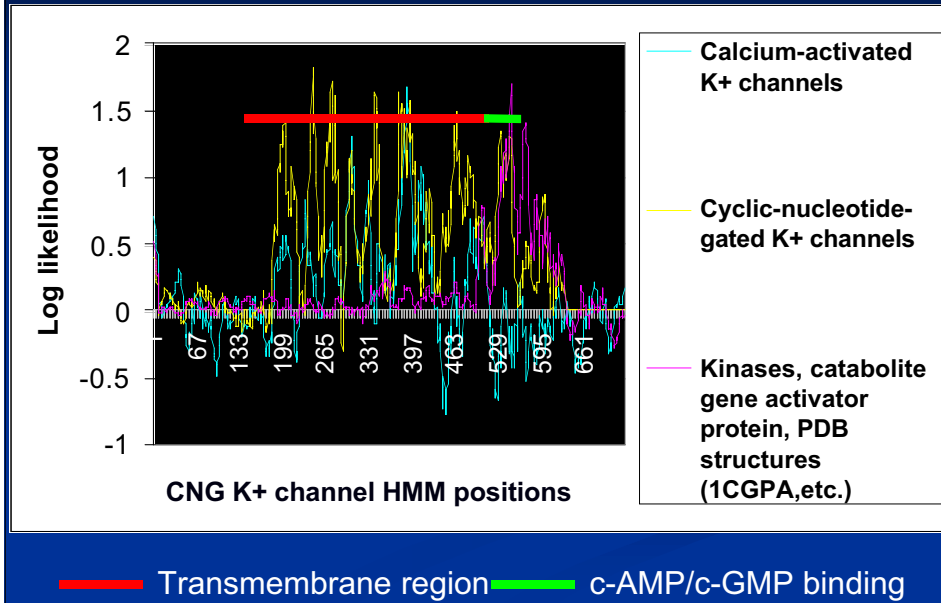**Cluster termination cutoff**

**Plot LL clusters**
**(simple average)**



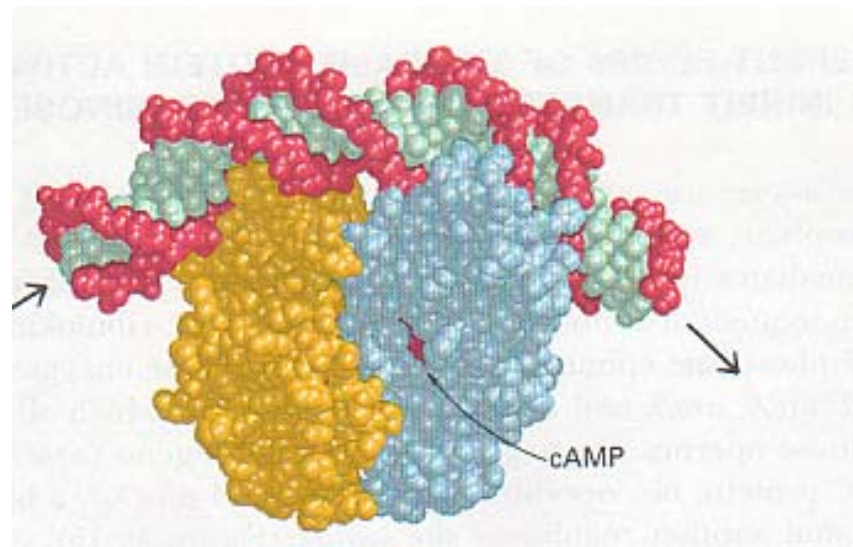Cyclic-Nucleotide-Gated K+ Channel
HMM Alignment Analysis

log likelihood

msa columns

Transmembrane region
c-AMP/c-GMP binding

MAG p22

## Analysis of CNG K+ domain structure



Log likelihood (y-axis: 2, 1.5, 1, 0.5, 0, -0.5, -1)

CNG K+ channel HMM positions (x-axis: 67, 133, 199, 265, 331, 397, 463, 529, 595, 661)

Legend:
- Calcium-activated K+ channels
- Cyclic-nucleotide-gated K+ channels
- Kinases, catabolite gene activator protein, PDB structures (1CGPA,etc.)

Transmembrane region  c-AMP/c-GMP binding

## Catabolite Gene Activator Protein (CAP) bound to DNA



cAMP

# Database hits to Voltage-gated K+ channel

| | |
|---|---|
| gi\|3023481\|sp\|CIKB_RAT | -427.39 |
| gi\|3913257\|sp\|CIKB_HUMAN | -426.47 |
| gi\|348462\|pir\|\|A44838 | -425.46 |
| gi\|345875\|pir\|\|S31761 | -423.14 |
| gi\|3023495\|sp\|CIKA_HUMAN | -422.38 |
| . . . | |
| gi\|1147595\|emb\|CAA64176.1\| | -32.56 |
| gi\|3874832\|emb\|CAA94204.1 | -32.50 |
| gi\|487428\|gb\|AAA50173.1\| | -32.49 |
| gi\|3452399\|gb\|AAC32857.1 | -32.41 |
| gi\|2648884\|gb\|AAB89577.1\| | -32.39 |
| gi\|2832781\|emb\|CAA12645.1 | -32.34 |
| gi\|1255396\|gb\|AAA96127.1\| | -32.13 |
| gi\|116452\|sp\|P15389\|CIN5_RAT | -32.09 |
| gi\|3924830\|emb\|CAA98957.1\| | -31.74 |
| gi\|465874\|sp\|P34410\|TWK8_CAEEL | -31.55 |
| gi\|2315751\|gb\|AAB66175.1\| | -31.48 |
| gi\|2665784\|gb\|AAC29515.1\| | -31.48 |
| gi\|2315752\|gb\|AAB66176.1\| | -31.45 |
| gi\|2315635\|gb\|AAB66084.1\| | -31.28 |
| gi\|1707203\|gb\|AAB37942.1\| | -31.23 |
| gi\|1181413\|gb\|AAC96618.1\| | -31.21 |
| gi\|3881291\|emb\|CAA21749.1\| | -31.14 |

**Ion channels**

**Similar to TNF-alpha-induced protein B12**

**Unknown**

?

## Analysis of Voltage-gated K+ channels domain structure

**TNF-alpha-induced protein B12 and K+ channel tetramerization domain (1T1DA)**



Legend:
— Tetramerization domain of K+ channels

— Similar to TNF-alpha-induced protein B12 (30 seqs)

x-axis: HMM positions
y-axis: log likelihood

## TNF-alpha acts on K+ current…but how?

# Predictions:

TNF-alpha-induced protein B12 and K+ channel tetramerization domain (1T1DA)

- Tetramerization domain of K+ channels
- Similar to TNF-alpha-induced protein B12 (30 seqs)
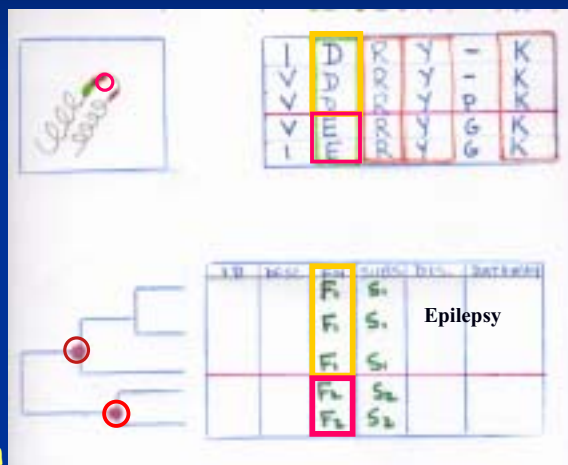
**1T1DA**
**Tetramerization Domain of K+ channels**

● Structure of TNF-alpha induced protein B12 is homologous to K+ channel tetramerization domain

● Does TNF-alpha induced protein B12 affect K+ channel function by interacting with the K+ channel T1 domain?

---

# Web server for high-throughput functional classification of proteins



3D structure viewer

Multiple sequence alignment

Epilepsy
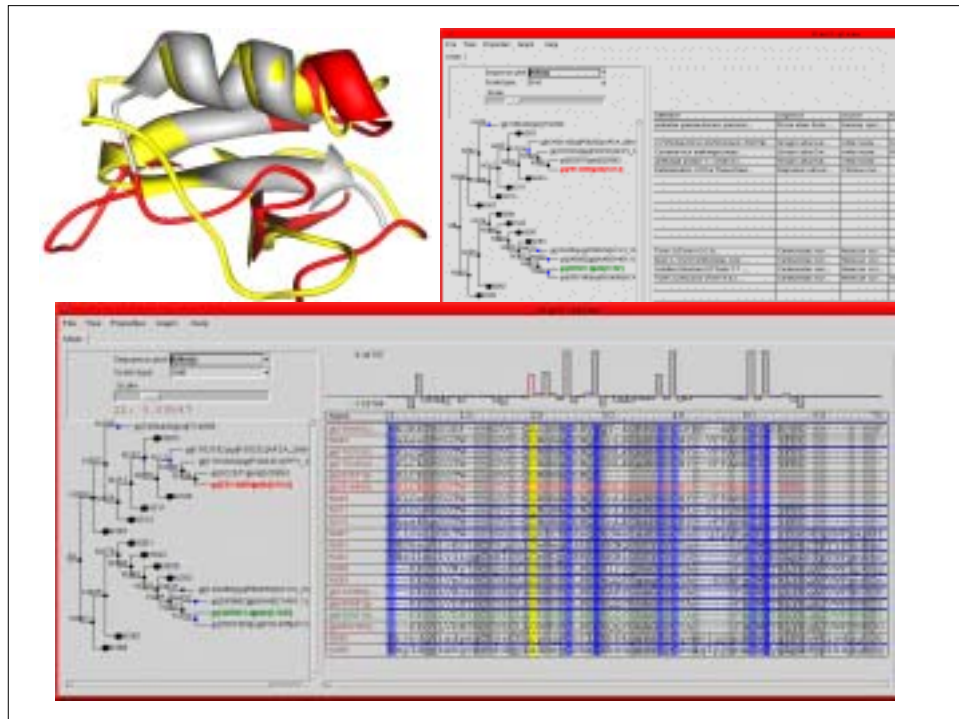
Phylogenetic tree, with subfamily decomposition

Attribute data table

Enable and foster virtual collaborations, scientific discovery, correction of errors in database annotations.
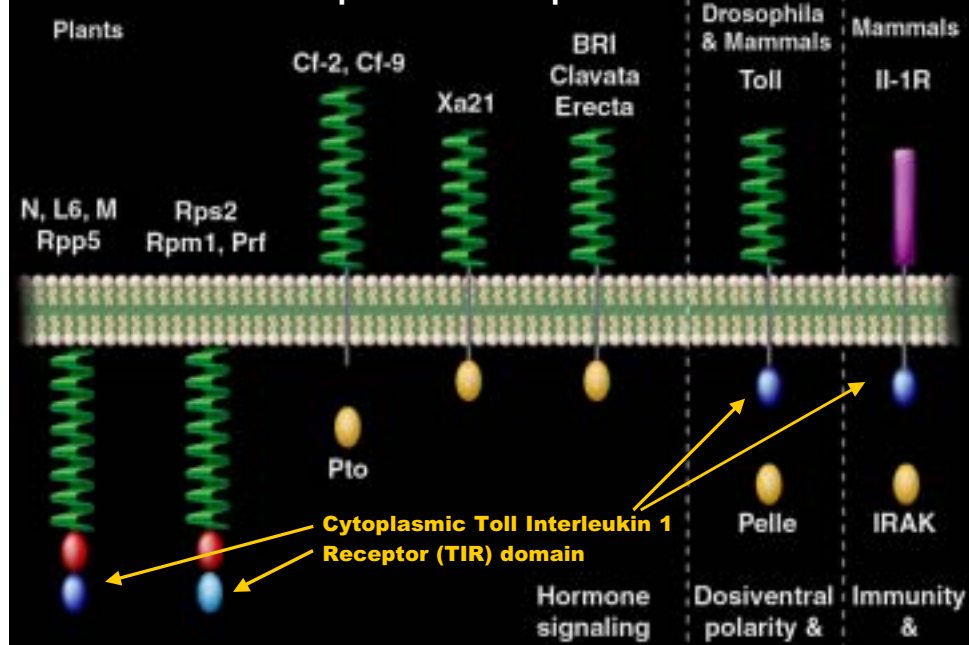
# Stalking disease-resistance proteins in rice

Toll Interleukin 1 domain
PDB structure 1FYX

Joint work with Barbara Baker and Brian Staskawicz

**Plant and Animal Innate Immunity Mediated by Structurally Similar Receptor and Receptor-like molecules**



**Conserved "scaffolding" proteins in cell death**

*C.elegans* CED- | CARD | NBS |

human Apaf-1 | CARD | NBS | WD40 |

human Nod1 | CARD | NBS | LRR |

tobacco N | TIR | NBS | LRR |

tobacco N$^{tr}$ | TIR | NBS |
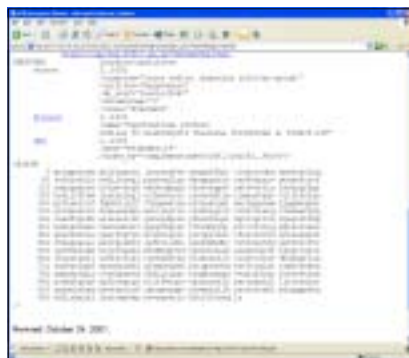
CARD: CAspase Recruitment Domain
Nod is implicated in CrohnÅs disease
N is involved in plant Hypersensitive Response (HR)

# TIR domains missing from monocot species…



ÅMonocot sequences are absent from the TIR subfamilyÅ
ÅToll/Interleukin-1 receptor homology (TIR), Å was entirely
absent from monocot species in searches of both random
genomic sequences and large collections of ESTsÅ

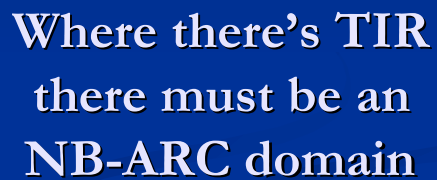# Searching the rice genome with general and subfamily HMMs for the TIR domain…



gi|14587298 gets top score (e-12 to TIR structure)

# Alignment of rice sequence to Toll-like receptor 2 (TLR2) subfamily

Alignment of gi|14587298|_Rice_putative_TIR_containing_protein with consensus sequence for subfamily

difya_ 4.23.2.1 (A) Toll-like receptor 2, TLR2 (Human (Homo sapiens))

N, L6, M
Rpp5

Rps2
Rpm1, Prf

LRR

NB-ARC

TIR

## Where there's TIR there must be an NB-ARC domain

## Alignment to APAF_Human (and homologs) NBS domain



## Alignment to HPr kinase (serine kinase/phosphatase) with P-loop/Walker motif A
### (possible structural similarity)

Phosphate binding motif

BLAST fails to detect TIR domain homologs for this sequence



Clustering and aligning homologs with FlowerPower

TIR domain

Walker A motif

NBS-ARC domain

# UC Berkeley Phylogenomics

**Group Members:**

David Konerding, Ph.D.

Wayne Christopher, Ph.D.

Bob Edgar, Ph.D.

Austin Huang

Joseph Dale

**Investigation of plant disease-resistance proteins**

Brian Staskawicz

Barbara Baker

Richard Michelmore

**Thanks to the National Science Foundation**

http://phylogenomics.berkeley.edu

kimmen@uclink.berkeley.edu